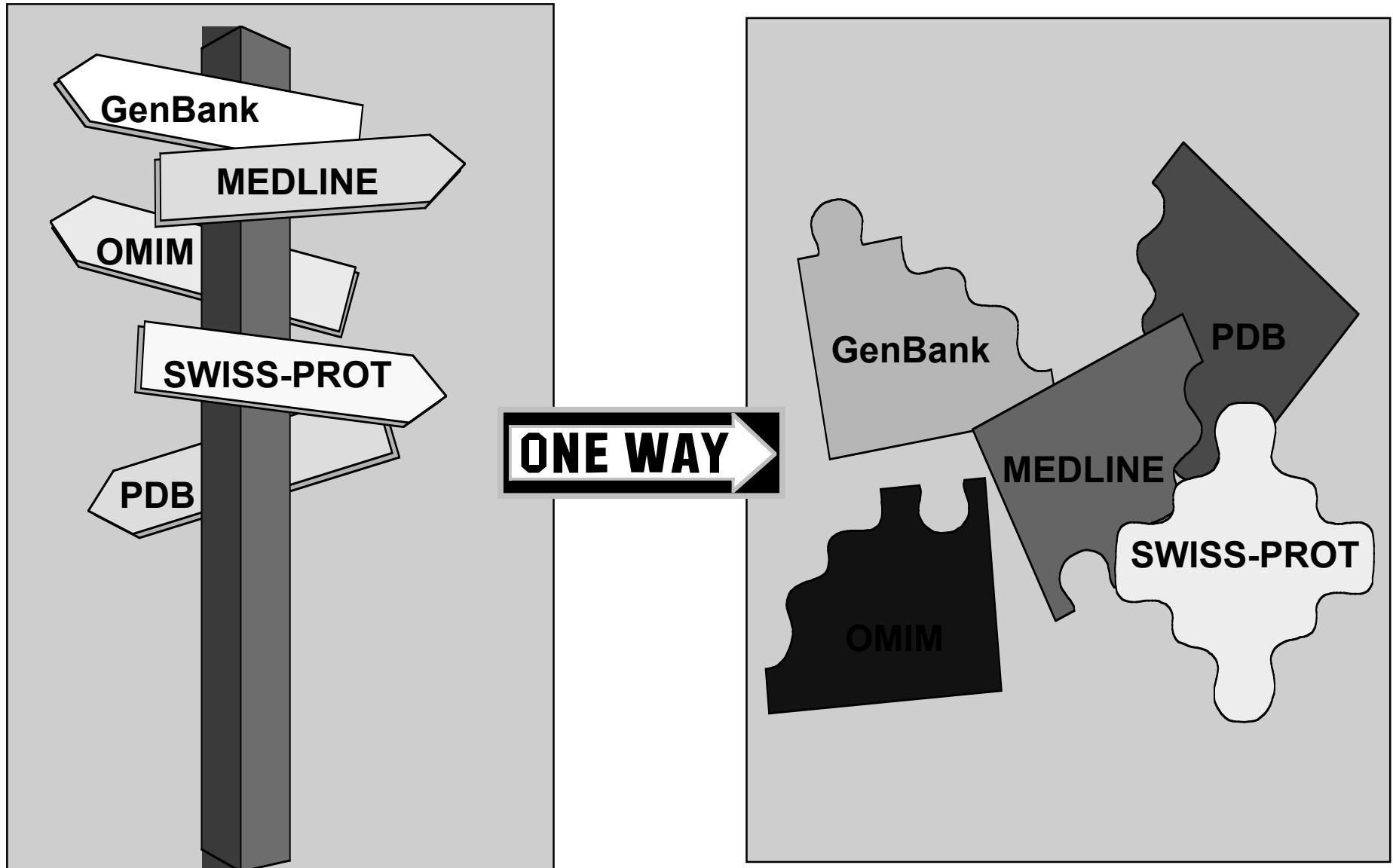
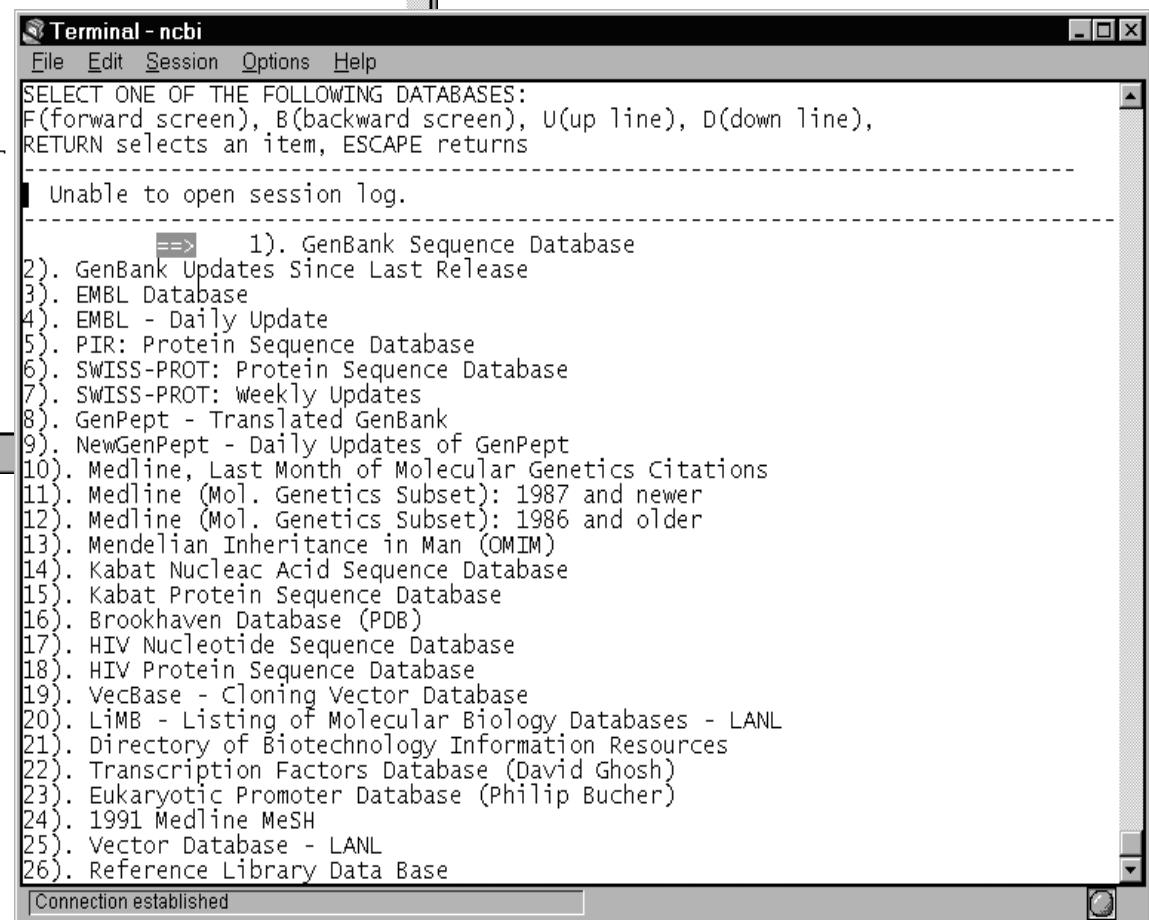
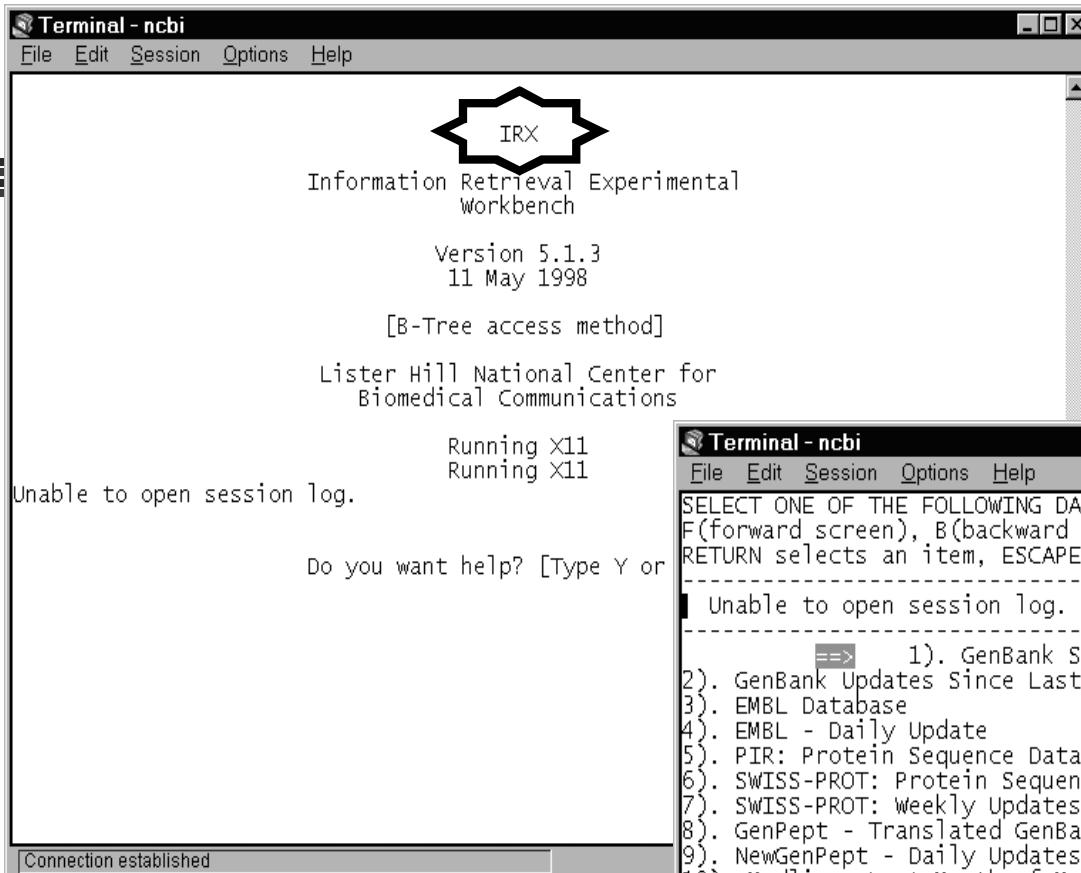

Biological Databases: Information Retrieval

David Landsman
Computational Biology Branch
NCBI

Information puzzles end abruptly?

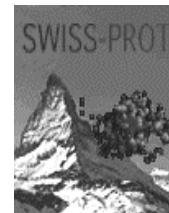




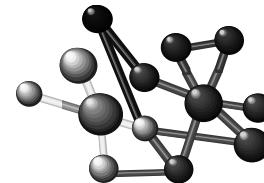
Vertical querying



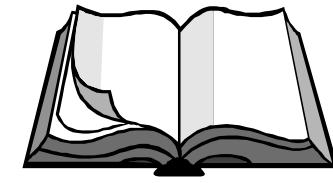
GenBank



SWISS-PROT



PDB

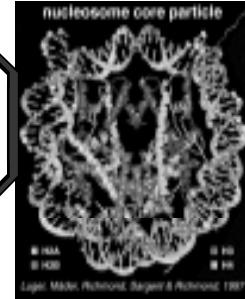


MEDLINE



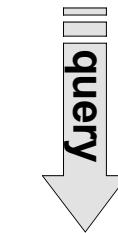
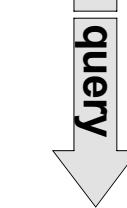
atg gcc cga acc aa
act gct cgt aag t
ggg aaa gcc
aaa cag ctg gcc
gcc gcc agg aaa agc

MARTKQTARI
TGGKAPRE
ATKAARKS

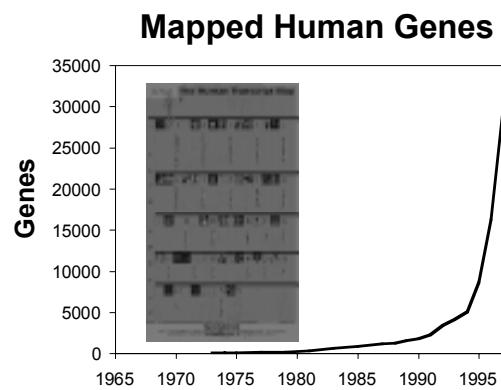
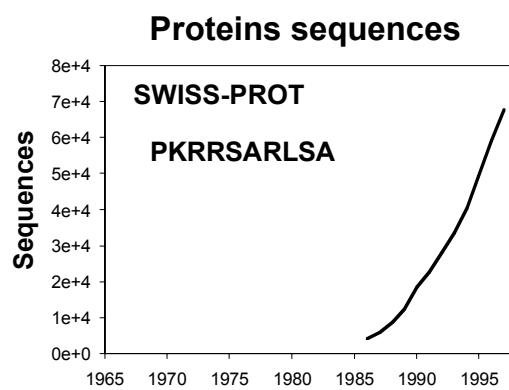
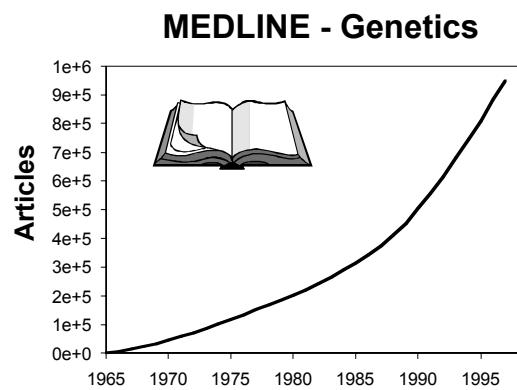
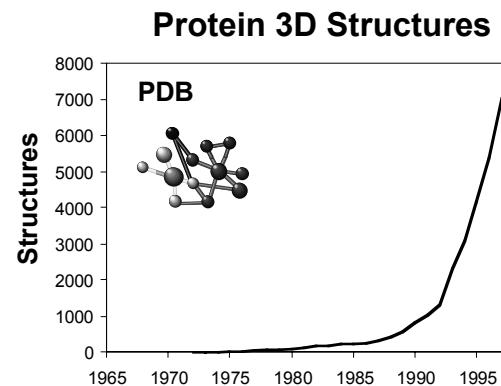
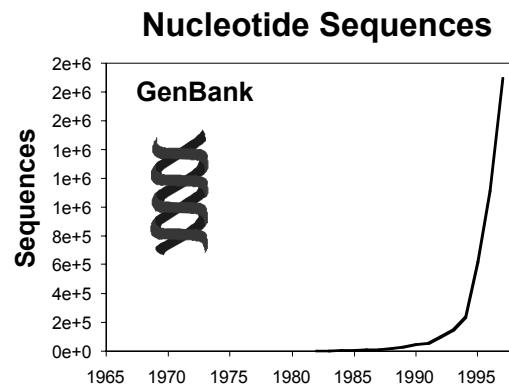
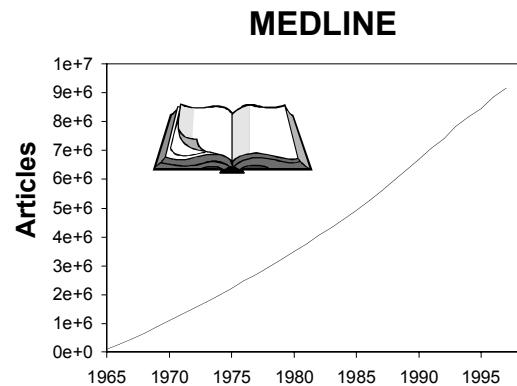


Nature
1997 Sep
18;389(6648):251-260
Crystal structure of the
nucleosome core particle
at 2.8 Å resolution.

Luger K, Mader AW,
Richmond RK, Sargent
DF,
Richmond TJ

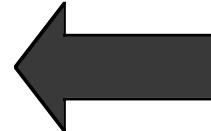


The Biotechnology Information Explosion



Entrez and its databases

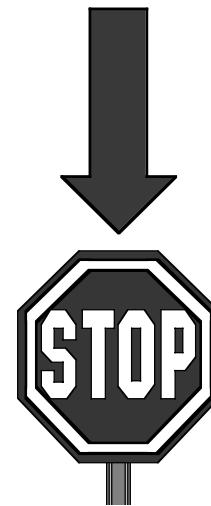
horizontal
or lateral
querying



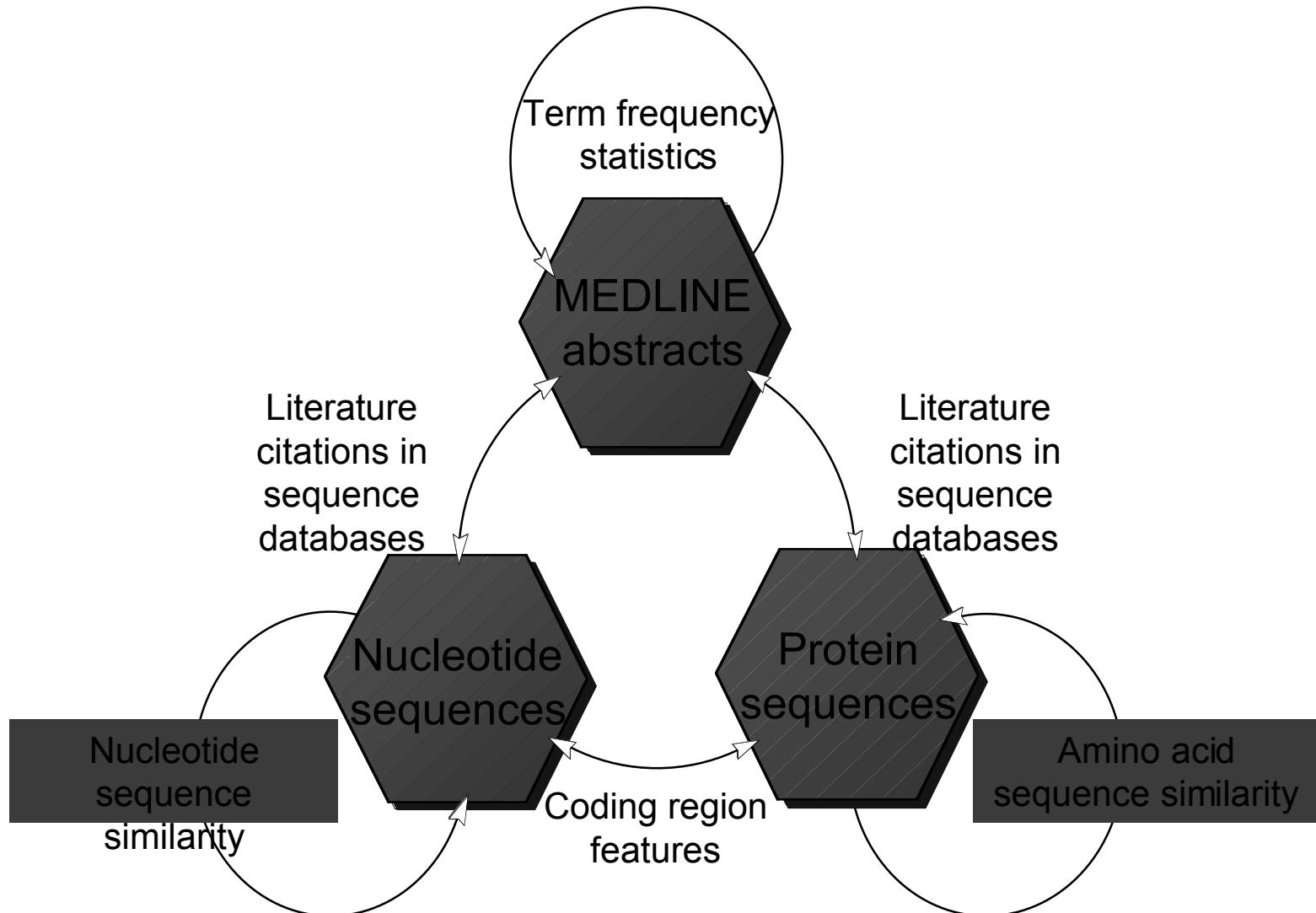
vertical
querying



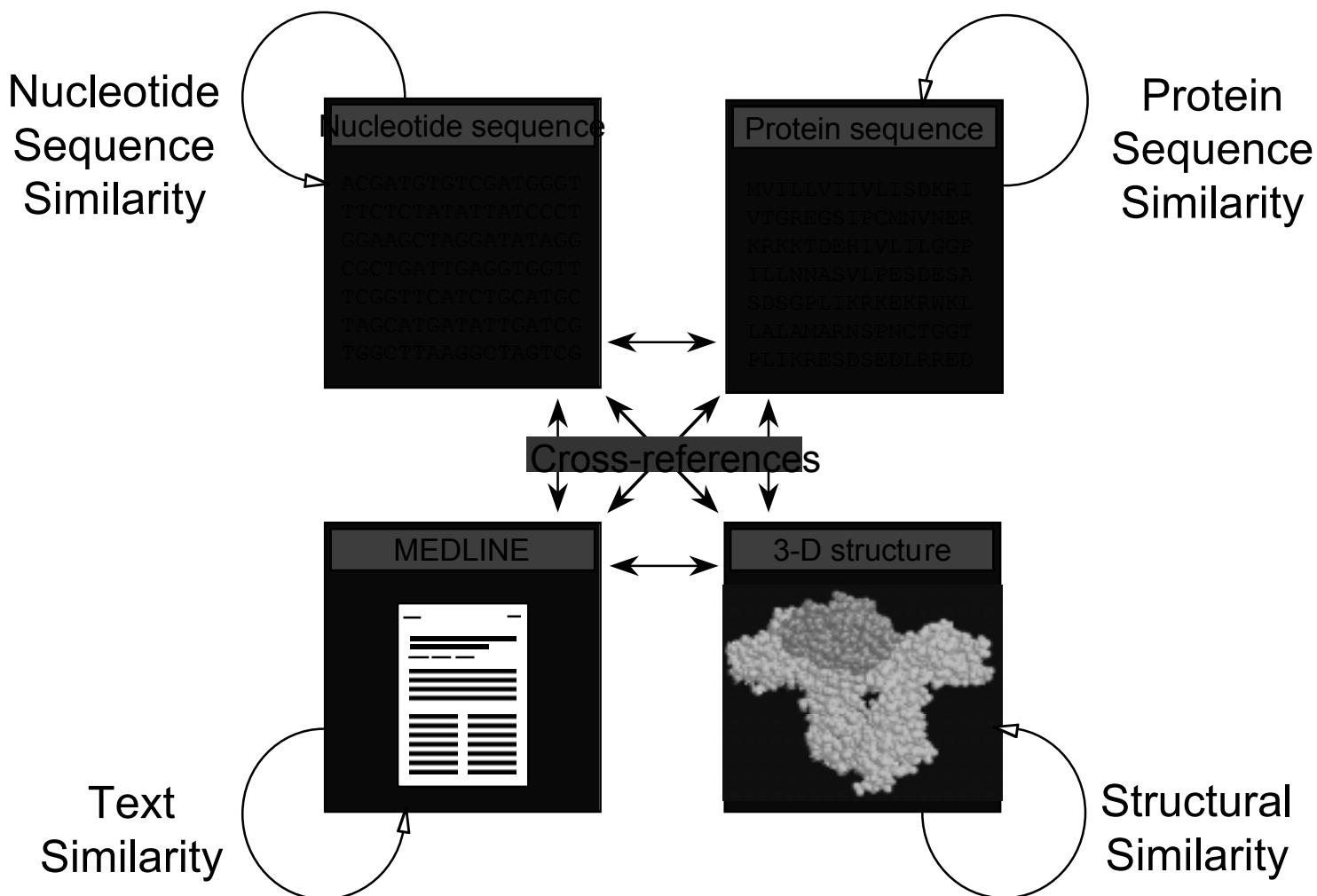
horizontal
or lateral
querying



Entrez (1992)



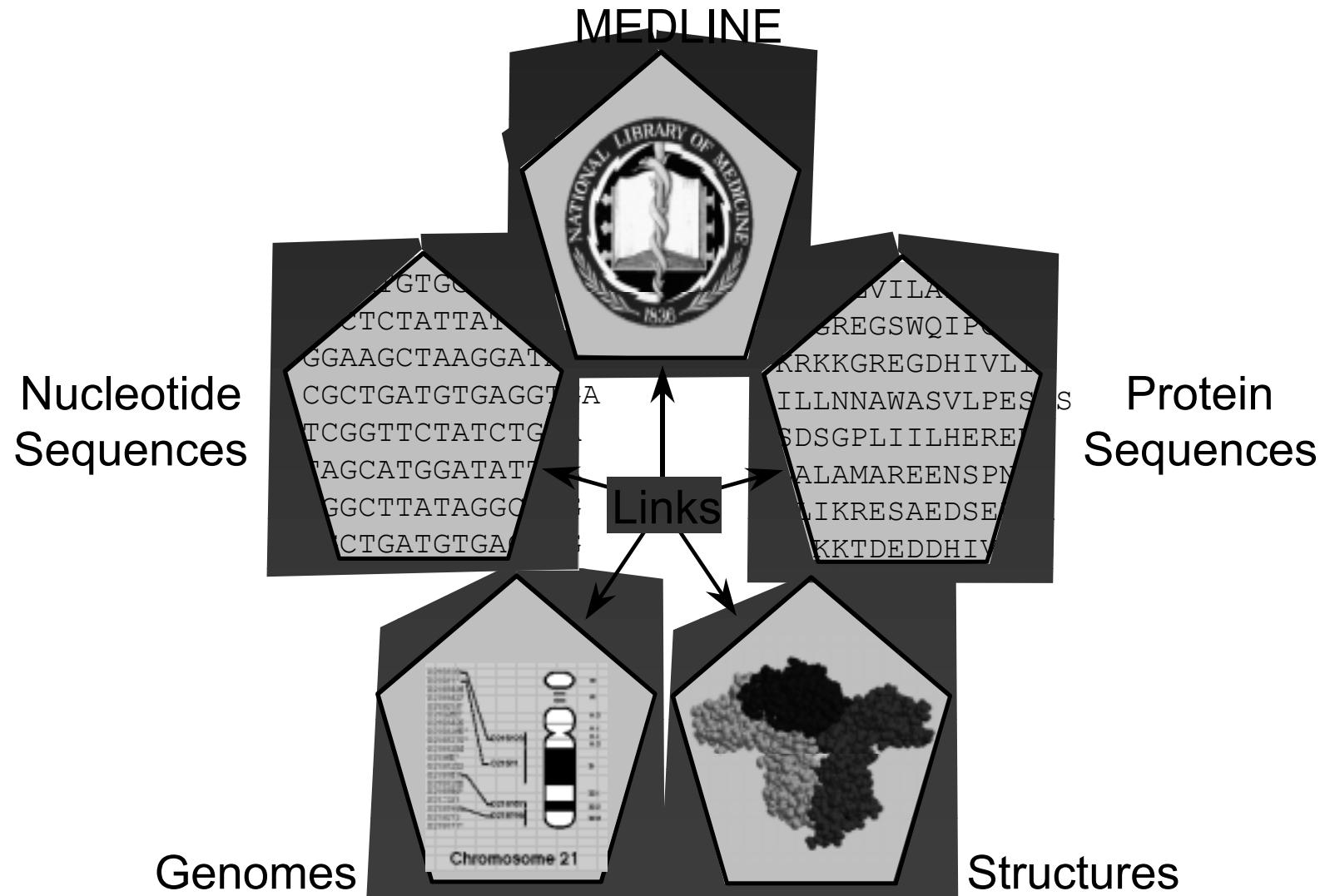
Entrez (1994)



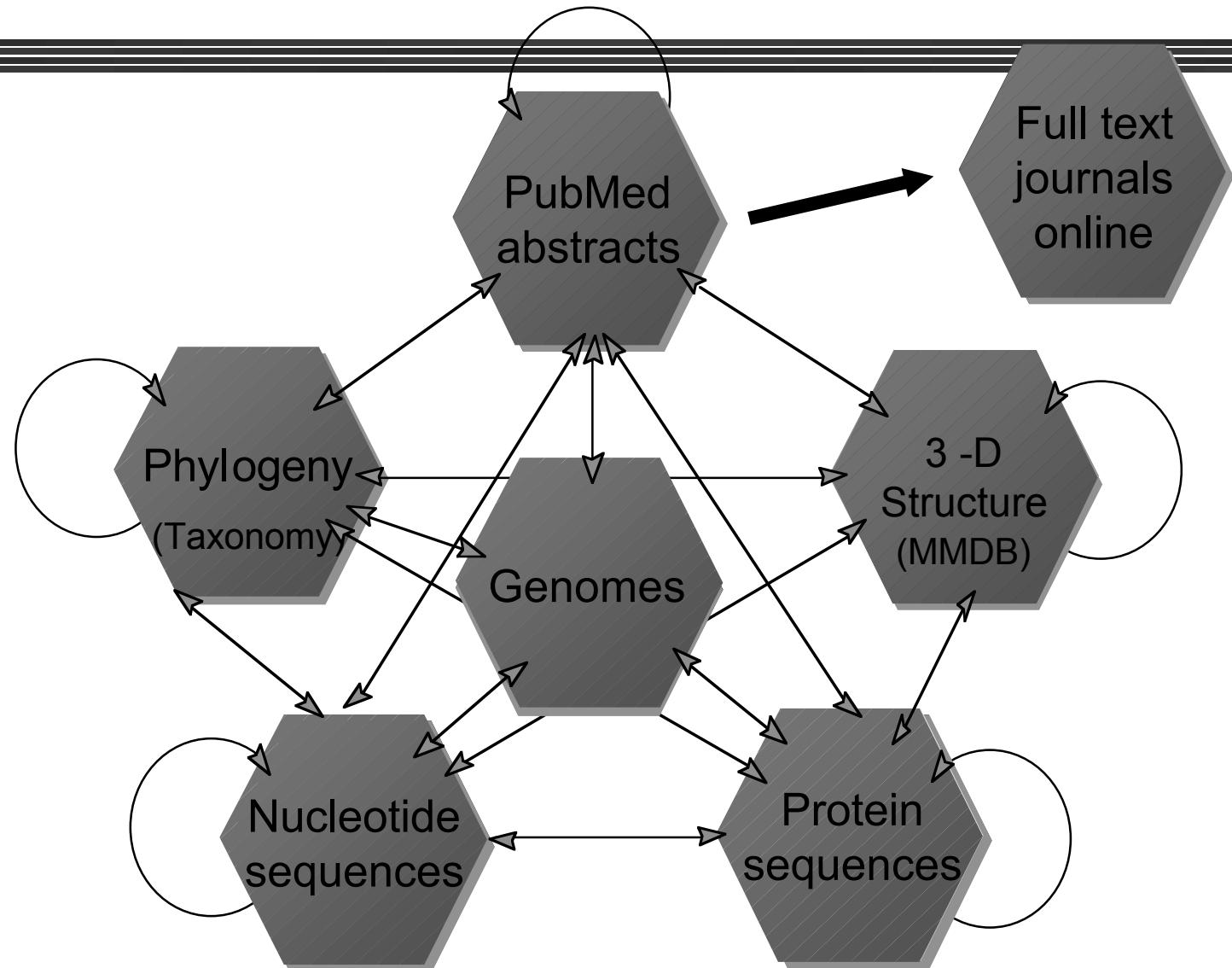
Entrez BigMed Data Set (Replaced by PubMed)

- Version 27.58; April 22, 1997
- 1,392,146 MEDLINE entries
 - All citations indexed under the MeSH term “Genetics” and its descendants (i.e. G5 tree)
 - Smaller subsets having to do with molecular sequences
 - Compare to total MEDLINE with ~ 8 million records
- 421,797 protein entries (“nr”)
- 1,273,697 nucleotide entries (“nr”)
- 4,873 structure entries (MMDB)
- 47,213 genome entries
- linked journals: *J. Biol. Chem.*, *PNAS*, *Science*, *J.M.B.*,
++++

Entrez: Genotype to Phenotype (1996)



Entrez Increases Discovery Space 1998



Entrez Implementations

■ CD-ROM subscription (stand-alone)

- Impractical data currency
 - ~100 new peer-reviewed articles per month
 - Sequence databases double ~20 months
- Slow access speed
- Capacity limited to 650 MB per CD-ROM

Entrez Implementations

■ Network Entrez (client-server)

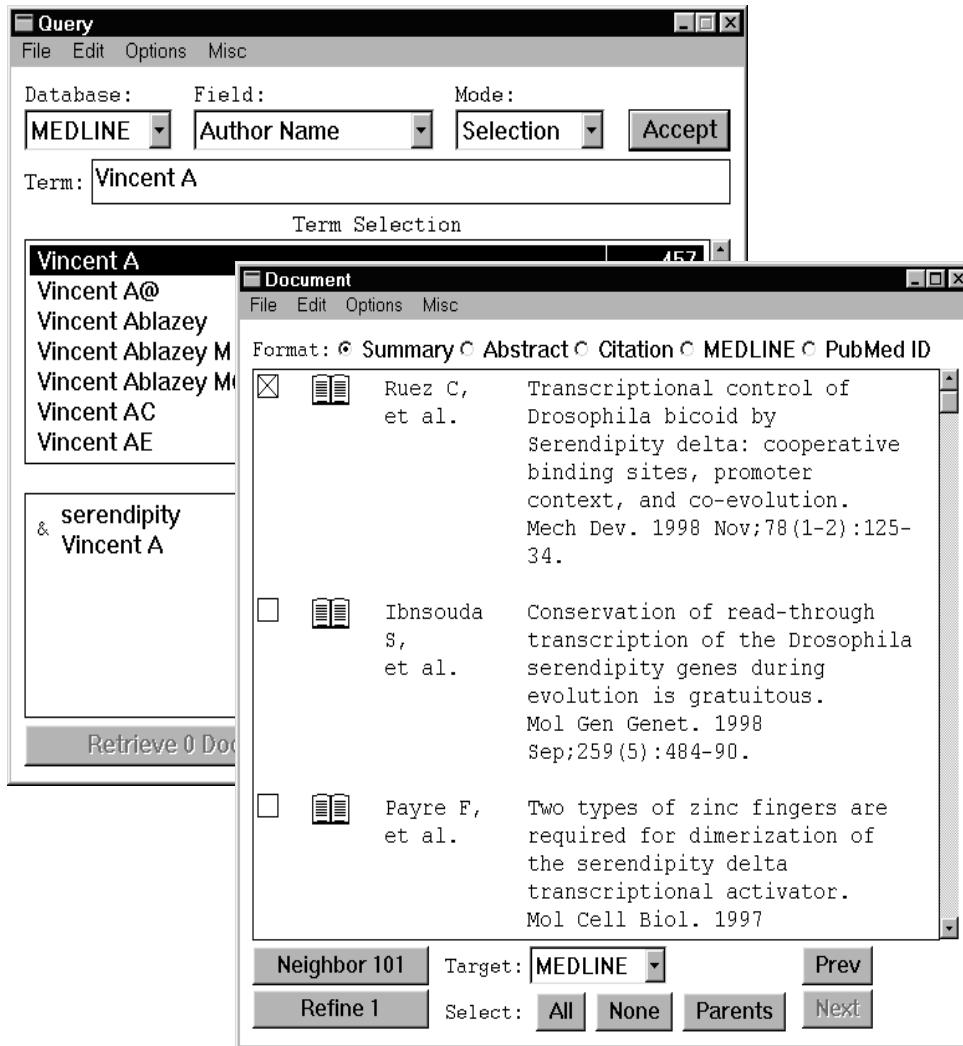
- Connect to dispatcher
- Fastest Entrez implementation
- Software installations and updates maintained by user

■ Web Entrez (WWW client-server)

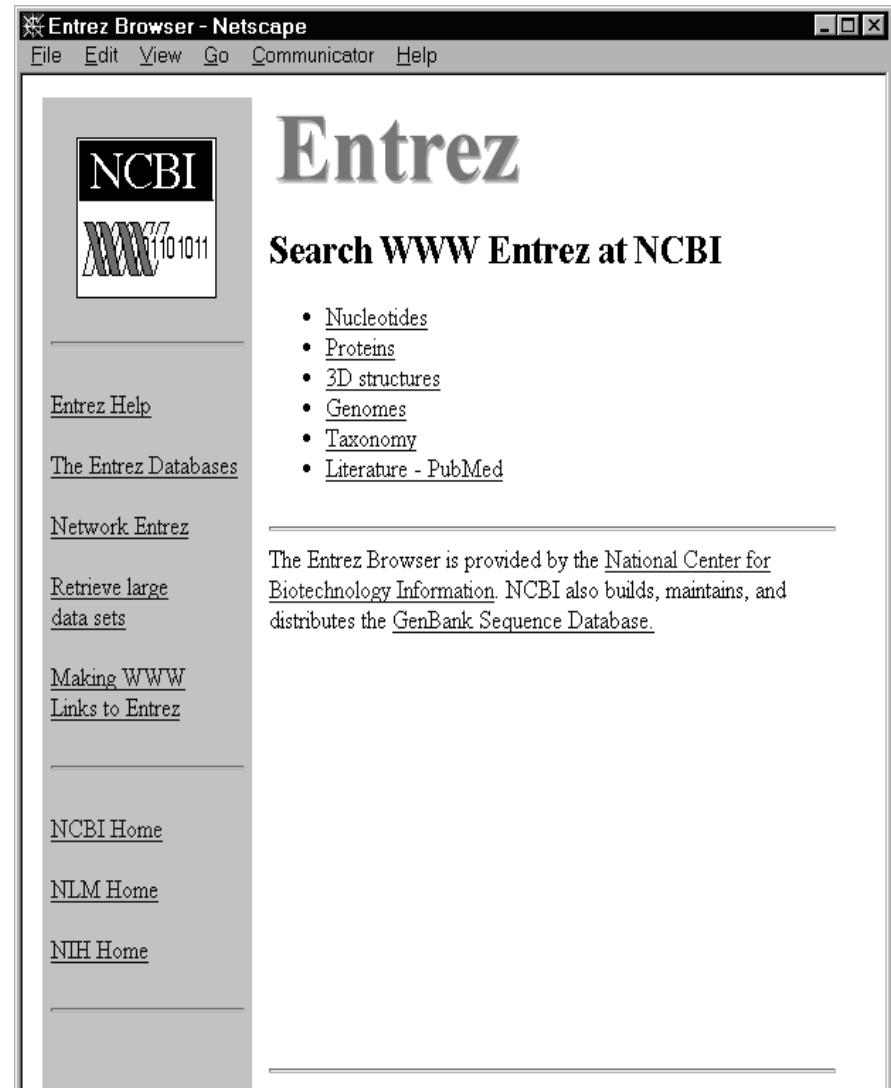
- Universality of WWW
- Can serve users without access to GUI environment (Lynx)
- Linking and neighboring information can be expressed as hypertext
- Ability to link to external data sources
- WWW browsers are supported by third-party developers
- Web implementation maintained by NCBI

Entrez Implementations

Network Entrez Version 6.6



Web Entrez



NCBI's Home Page

The screenshot shows the NCBI homepage as it appeared in Netscape Communicator. At the top, the title 'The National Center for Biotechnology Information - Netscape' is displayed above a menu bar with options: File, Edit, View, Go, Communicator, Help. Below the menu is the NCBI logo, which consists of a stylized 'X' made of horizontal bars followed by the text '01101011'. To the right of the logo is the text 'National Center for Biotechnology Information' and 'National Library of Medicine / National Institutes of Health'. A horizontal navigation bar below the title includes links for PubMed, Entrez, BLAST, BankIt, OMIM, Taxonomy, and Structure. Four large black arrows point upwards from the bottom of the page towards these menu items. The main content area is divided into several sections: 'Welcome to NCBI' (with links to What's New, Programs and Activities, NCBI Newsletter, NCBI Exhibit Schedule, Service Addresses, and User Comments); 'GenBank Sequence Database' (with links to Overview, Searching GenBank, Submitting Sequences, BankIt, and Sequin); 'Database Services' (with links to PubMed MEDLINE, Entrez Search System, BLAST Sequence Similarity Searching, E-mail Servers, and Anonymous FTP); 'Research' (with links to Research Projects, Seminar Schedule, Staff Bibliography, and Full Text of Selected Staff Publications); and 'Other NCBI Resources' (with a long list of links including Cancer Genome Anatomy Project, Gene Map of the Human Genome, Genes and Disease, UniGene, Human/Mouse Homology Maps, HGSI, Clusters of Orthologous Groups, OMIM, dbEST, dbGSS, dbSTS, dbSNP, Electronic PCR, MMDB, NCBI Taxonomy, ORF Finder, Malaria Genetics and Genomics, and Retrovirus Resources). At the bottom left is a 'Revised March 11, 1999' note. On the right side, a large rectangular box contains the URL 'www.ncbi.nlm.nih.gov'.

The National Center for
Biotechnology Information

National Library of Medicine / National Institutes of Health

PubMed Entrez BLAST BankIt OMIM Taxonomy Structure

Welcome to NCBI

- [What's New](#)
 - [Site Search](#)
 - [Organism-specific BLAST](#)
- [Programs and Activities](#)
- [NCBI Newsletter](#)
- [NCBI Exhibit Schedule](#)
- [Service Addresses](#)
- [User Comments](#)

GenBank Sequence Database

- [Overview](#)
- [Searching GenBank](#)
- [Submitting Sequences](#)
 - [BankIt](#)
 - [Sequin](#)

Database Services

- [PubMed MEDLINE](#)
- [Entrez Search System](#)
- [BLAST Sequence Similarity Searching](#)
- [E-mail Servers](#)
- [Anonymous FTP](#)

Research

- [Research Projects](#)
- [Seminar Schedule](#)
- [Staff Bibliography](#)
- [Full Text of Selected Staff Publications](#)

Other NCBI Resources

- [Cancer Genome Anatomy Project](#)
- [Gene Map of the Human Genome](#)
- [Genes and Disease](#)
- [UniGene: Unique Gene Sequence Collection for Human, Mouse, and Rat](#)
- [Human/Mouse Homology Maps](#)
- [HGSI: Human Genome Sequencing Index](#)
- [Clusters of Orthologous Groups](#)
- [OMIM: Online Mendelian Inheritance in Man](#)
- [dbEST: Database of Expressed Sequence Tags](#)
- [dbGSS: Database of Genome Survey Sequences](#)
- [dbSTS: Database of Sequence Tagged Sites](#)
- [dbSNP: Database of Single Nucleotide Polymorphisms](#)
- [Electronic PCR](#)
- [MMDB: Molecular Modelling Database](#)
- [NCBI Taxonomy](#)
- [ORF Finder](#)
- [Malaria Genetics and Genomics](#)
- [Retrovirus Resources](#)

Revised March 11, 1999

www.ncbi.nlm.nih.gov

Advanced PubMed querying

<http://www.ncbi.nlm.nih.gov/PubMed/medline.html>

Advanced PubMed Search - Netscape

File Edit View Go Communicator Help

NATIONAL LIBRARY OF MEDICINE

PubMed

Advanced MEDLINE Search

Search Field: All Fields Mode: Automatic

Affiliation

All Fields

Author Name
EC/RN Number
Entrez Date
Issue
Journal Name
Language
Journal MeSH Major Topic
• Enter o
• Author
• Journal
• Journal MeSH Terms
• Use pu
• Boolean
NOT 1
Page
Publication Date
Publication Type
Subheading
Substance Name
Text Word
Title Word
Volume
MEDLINE ID
PubMed ID

Search

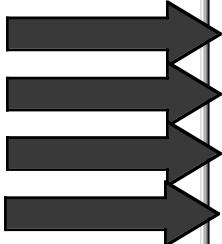
• Enter o
• Author
• Journal
• Journal MeSH Terms
• Use pu
• Boolean
NOT 1
Page
Publication Date
Publication Type
Subheading
Substance Name
Text Word
Title Word
Volume
MEDLINE ID
PubMed ID

is.
ed in the form Smith JB, but initials are optional.
n full, as valid MEDLINE abbreviations, or as ISSN numbers (see
rmation).
fy fields and search mode.
also be entered directly in the search box, using AND, OR, and
Boolean Search page for more information.

Questions or comments? [Help Desk](#).

NOTICE: The PubMed data are for personal use only. Users are responsible for complying with all copyright and licensing restrictions associated with data.

Overview
Help / FAQs
New/Noteworthy
Clinical Alerts
Basic Search
Clinical Queries
Journal Browser NEW
MeSH Browser
Citation Matcher
Internet Grateful Med
NLM Home
NCBI Home
NIH Home
Credits
Restrictions on use



PubMed Clinical Queries - Netscape

File Edit View Go Communicator Help

MeSH Browser - Netscape

File Edit View Go Communicator Help

NCBI PubMed MeSH Browser

PubMed ?

Enter another MeSH term

Histones [Detailed]

Small chromosomal DNA in cell nuclei by is based on the relative arrangement of nucleosomes.

Note! The term 'Histone' is currently in the MeSH tree at location 2.

Jump to: [location 2](#)

MeSH Tree Location 1: Histones

Top of MeSH Tree

- [Amino Acids, Peptides, Proteins](#)
- [Nuclear Proteins](#)
- [Histones](#)

MeSH Tree Location 2: Histones

Top of MeSH Tree

- [Amino Acids, Peptides, Proteins](#)
- [Nucleoproteins](#)
- [Histones](#)

Journal Database Browser - Netscape

File Edit View Go Communicator Help

NCBI Journal Database Browser

Entries ?

Found 123 journal(s)

Prev matches

Title
PROGRESS IN FOOD AND NUTRITION SCIENCE
PROTEIN SCIENCE
QUARTERLY JOURNAL OF MICROSCOPICAL SCIENCE
RESEARCH IN VETERINARY SCIENCE
ROGERSIAN NURSING SCIENCE NEWS
SCANDINAVIAN JOURNAL OF MEDICINE AND SCIENCE
SCIENCE AND JUSTICE
SCIENCE ET RECHERCHE ODONTOSTOMATOLOGIQUE
SCIENCE IN CHINA. SERIES B, CHEMISTRY, LIFE SCIENCES

Next matches

Enter the journal name, MEDLINE abbreviation or ISSN. You may enter up to 10 words in the title. Use an asterisk (*) at the end of a word to tell us you want all words starting with that stem.

All types Science

[Start...](#)

Comments and questions to the [Help Desk](#)
Credit: Andrey Smirnov

Citation Matcher for Single Articles - Netscape

File Edit View Go Communicator Help

PubMed Citation Matcher for Single Articles

Enter information about the article you wish to find.

Journal

Date:

Volume: Issue: First page:

Author's last name and initials (e.g., Smith BJ)

[Search](#) [Clear](#)

Notes:

- You may omit any item if you wish.
- Journal titles may be entered in full or as valid MEDLINE abbreviations.
- For Date, you may enter yyyy, yyyy/mm, or yyyy/mm/dd. For example, 1998, 1998/03, or 1998/03/06.
- Author names are automatically truncated to account for varying initials, e.g., smith j will also match on smith ja, smith jb, smith jc jr, etc. Enclose author names in double quotes to retrieve that exact match, e.g., "smith j"

Comments and questions to the [Help Desk](#)

Entrez Genomes Division

The screenshot shows a vintage web browser window titled "Entrez Browser - Netscape". The menu bar includes "File", "Edit", "View", "Go", "Communicator", and "Help". The main content area features the NCBI logo and the word "Entrez" in large, serif capital letters. Below it is the heading "Search WWW Entrez at NCBI" and a bulleted list of search categories: "Nucleotides", "Proteins", "3D structures", "Genomes", "Taxonomy", and "Literature - PubMed". To the left of the main content is a sidebar with links: "Entrez Help", "The Entrez Databases", "Network Entrez", "Batch Entrez (Retrieve large data sets)", and "Making WWW Links to Entrez". At the bottom of the sidebar are links to "NCBI Home", "NLM Home", and "NIH Home". A footer at the bottom of the page says "Comments and questions to the [NCBI help desk](#)".

Entrez genomes - Netscape

File Edit View Go Communicator Help

NCBI Entrez Genomes

NCBI BLAST Nucleotides Proteins Structure Taxonomy PubMed Help

Search_for: All Fields

All Organisms

Prominent Organisms

Microbial genomes

- BLAST
- List of projects

Archaea

- Genome
- Plasmids

Bacteria

- Genome
- Plasmids

Eukaryota

- Genome
- Plasmids
- Organelles

Viruses

Viroids

Plasmids

Prominent Organisms Taxonomy / List

Complete Genomes

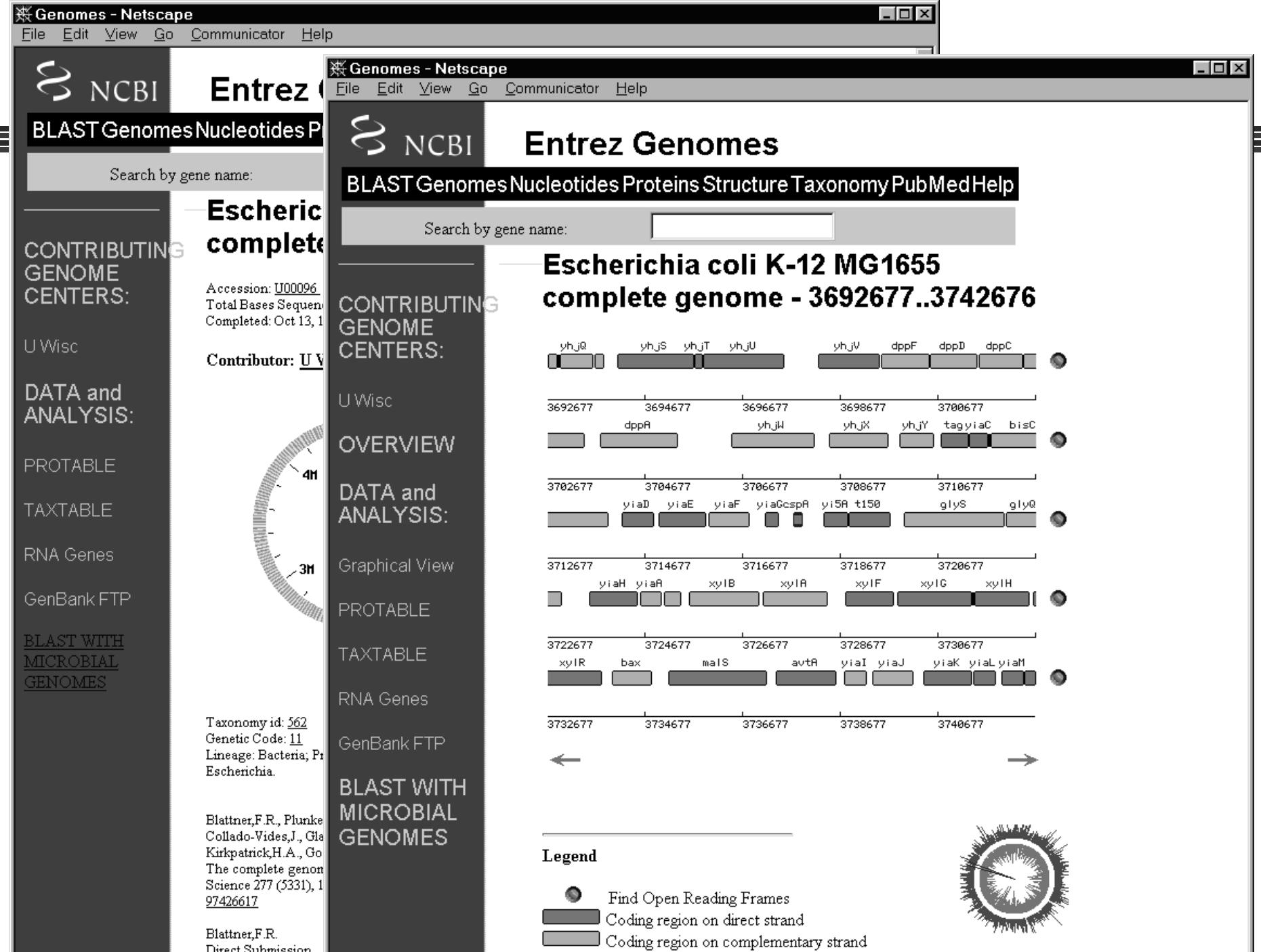
- *Aquifex aeolicus*
- *Archaeoglobus fulgidus*
- *Bacillus subtilis*
- *Borrelia burgdorferi*
 - *chromosome*,
 - plasmids: *cp26*, *cp9*, *lp17*, *lp25*, *lp28-1*, *lp28-2*, *lp28-3*, *lp28-4*, *lp36*, *lp38*, *lp54*.
- *Chlamydia trachomatis*
- *Chlamydia pneumoniae* NEW
- *Escherichia coli*
- *Haemophilus influenzae*
- *Helicobacter pylori*
- *Helicobacter pylori J99*
- *Methanobacterium thermoautotrophicum*
- [3] *Methanococcus jannaschii*
 - *chromosome*
 - *small extrachromosomal element*
 - *large extrachromosomal element*
- *Mycobacterium tuberculosis*
- *Mycoplasma genitalium*
- *Mycoplasma pneumoniae*
- *Pyrococcus horikoshii*
- *Rickettsia prowazekii*
- [16] *Saccharomyces cerevisiae*
 - chromosomes: I, II, III, IV, V, VI, VII, VIII, VIII, IX, X, XI, XII, XIII, XIV, XV, XVI.
- *Synechocystis PCC6803*
- *Treponema pallidum*

Complete Sequences

- *Leishmania major chromosome 1*
- *Plasmodium falciparum chromosome 2*

Chromosome Maps

- [6] *Caenorhabditis elegans*
 - chromosomes: I, II, III, IV, V, X.
- [5] *Drosophila melanogaster*
 - chromosomes: 1, 2, 3, 4, Y.
- [23] *Homo sapiens*
 - chromosomes: 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, X.
- [21] *Mus musculus*



Site Name - Netscape

File Edit View Go Communicator Help

NCBI NCBI Entrez Genomes BLAST Taxonomy Structure

Search by gene name: HHO1

Saccharomyces cerevisiae complete genome

	100 300 500 700 900 1100 1300 1500 1700 1900	
I	[chromosome diagram]	229,237
II	[chromosome diagram]	813,138
III	[chromosome diagram]	315,339
IV	[chromosome diagram]	1,531,974
V	[chromosome diagram]	576,870
VI	[chromosome diagram]	270,148
VII	[chromosome diagram]	
VIII	[chromosome diagram]	
IX	[chromosome diagram]	
X	[chromosome diagram]	
XI	[chromosome diagram]	
XII	[chromosome diagram]	
XIII	[chromosome diagram]	
XIV	[chromosome diagram]	
XV	[chromosome diagram]	
XVI	[chromosome diagram]	

FTP site

MIPS

SGD

YPD

Site Name - Netscape

File Edit View Go Communicator Help

NCBI NCBI Entrez Genomes BLAST Taxonomy Structure

Search by gene name:

Caenorhabditis elegans complete genome

	1 3 5 7 9 11 13 15 17 19 21	
I	[chromosome diagram]	12,520,182
II	[chromosome diagram]	16,264,216
III	[chromosome diagram]	11,417,571
IV	[chromosome diagram]	10,730,331
V	[chromosome diagram]	20,625,436
X	[chromosome diagram]	16,009,602

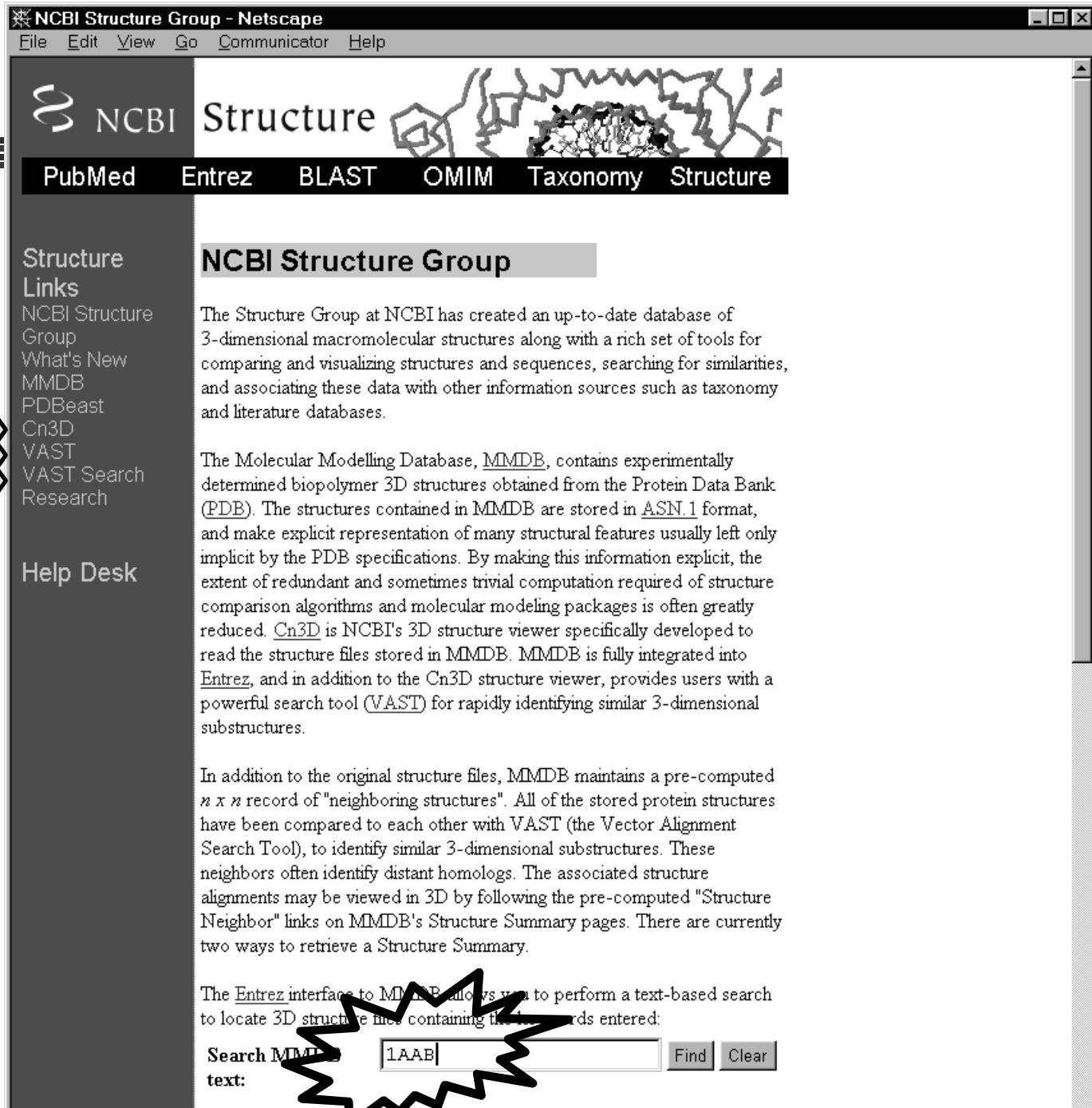
FTP SITE

Genome Sequencing:

Sanger Centre WashU

Structure searching ...

<http://www.ncbi.nlm.nih.gov/Structure/>



The Structure Group at NCBI has created an up-to-date database of 3-dimensional macromolecular structures along with a rich set of tools for comparing and visualizing structures and sequences, searching for similarities, and associating these data with other information sources such as taxonomy and literature databases.

The Molecular Modelling Database, [MMDB](#), contains experimentally determined biopolymer 3D structures obtained from the Protein Data Bank ([PDB](#)). The structures contained in MMDB are stored in [ASN.1](#) format, and make explicit representation of many structural features usually left only implicit by the PDB specifications. By making this information explicit, the extent of redundant and sometimes trivial computation required of structure comparison algorithms and molecular modeling packages is often greatly reduced. [Cn3D](#) is NCBI's 3D structure viewer specifically developed to read the structure files stored in MMDB. MMDB is fully integrated into [Entrez](#), and in addition to the Cn3D structure viewer, provides users with a powerful search tool ([VAST](#)) for rapidly identifying similar 3-dimensional substructures.

In addition to the original structure files, MMDB maintains a pre-computed $n \times n$ record of "neighboring structures". All of the stored protein structures have been compared to each other with VAST (the Vector Alignment Search Tool), to identify similar 3-dimensional substructures. These neighbors often identify distant homologs. The associated structure alignments may be viewed in 3D by following the pre-computed "Structure Neighbor" links on MMDB's Structure Summary pages. There are currently two ways to retrieve a Structure Summary.

The [Entrez](#) interface to MMDB allows you to perform a text-based search to locate 3D structure files containing the keywords entered:

Search MMDB text:

Cn3D Home Page - Netscape

File Edit View Go Communicator Help

NCBI Structure

PubMed Entrez BLAST OMIM Taxonomy Structure

Related Links

- Cn3D Tutorial
- Cn3D Help
- Cn3D FAQ
- Cn3D Install

Structure Links

- NCBI Structure Group
- What's New
- MMDB
- PDBeast
- Cn3D
- VAST
- Research

Help Desk

Cn3D 2.0

Cn3D is a helper application for your web browser that allows you to view three dimensional structures from NCBI Entrez retrieval service. Cn3D runs on Windows, Mac OS, and Unix. The new version of Cn3D simultaneously displays structure, sequence, and alignment:

Cn3D 2.0

File Edit View Structure Color Help

1DOI	1AWD	10	20
		PTVEYLNlyev	vddngwdmyd
	1	YKVTLKTPsg	
		30	40
1DOI	1AWD	21	ddvfgreasdm dlddedYGSL
		11	EETI
		50	60

In the above example, two ferredoxins, 1DOI and 1AWD, are aligned.

Document Done

MMDB Structure Summary - Netscape

File Edit View Go Communicator Help

NCBI MMDB STRUCTURE SUMMARY Entrez ?

MMDB Id: 4475 PDB Id: 1AAB

Protein Chains: (single chain)

MEDLINE: PubMed

Taxonomy: Rattus norvegicus

PDB Authors: C.H.Hardman, R.W.Broadhurst, A.R.C.Raine, K.D.Grasser, J.O.Thomas & E.D.Laue

PDB Deposition: 28-Oct-95

PDB Class: Dna-Binding

PDB Title: Nmr Structure Of Rat Hmg1 Hmga Fragment

Sequence Neighbors: (single chain)

Structure Neighbors: (single chain)

View / Save Structure NEW [Get Cn3D 2.0 Now!](#)

Options: Viewer: Complexity:

Launch Viewer Cn3D v2.0 (asn.1) Cn3D Subset Up to 5 Models
 See File Cn3D v1.0 (asn.1) Virtual Bond Model Up to 10 Models
 Save File Mage All Atom Model All Models
 RasMol (PDB)

Help [MMDB](#), [Cn3D](#), [Viewing Options](#), [VAST](#), [PDB](#), [NCBI Structure](#)

Vast Results - Netscape

File Edit View Go Communicator Help

NCBI VAST STRUCTURE NEIGHBORS Entrez ?

Structures similar to MMDB 4475, 1AAB

Nmr Structure Of Rat Hmg1 Hmga Fragment

View / Save Alignments NEW [Get Cn3D 2.0 Now!](#)

Options: Launch Viewer See File Save File

Viewer: Cn3D v2.0 (asn.1) Mage (Kinemage) (PDB)

Complexity: Aligned Chains only Alpha Carbons only
 All Chains All Atoms

Structure neighbors 1-3 out of 3 displayed. Page 1 of 1.

	PDB	C	D	RMSD	NRES	%Id	Description
<input type="checkbox"/>	1HJO	A	3	2.6	43	14.0	Heat-Shock 70kd Protein 42kd Atpase N-Terminal Domain
<input type="checkbox"/>	1A81	C	5	1.8	26	3.8	Crystal Structure Of The Tandem Sh2 Domain Of The Syk Kinase Bound To A Dually Tyrosine-Phosphorylated Itam
<input checked="" type="checkbox"/>	2LEF	A		1.7	55	12.7	Lef1 Hmg Domain (From Mouse), Complexed With Dna (15bp), Nmr, 12 Structures

Display / Sort Hits page number: Hits to display per page: choose between 20-100 neighbors per page.

Display Subset:

- Non-redundant; BLAST p-value 10e-7
- Non-redundant; BLAST p-value 10e-40
- Non-redundant; BLAST p-value 10e-80
- Non-identical sequences
- All of MMDB

Sorted by:

- VAST Score
- VAST P-value
- Rmsd
- Aligned residues
- Identities

Column Format:

- RMSD, NRES, %Id
- All values



